

WHAT IS CLAIMED IS:

1. A communications system, comprising:

(a) a plurality of switch nodes, each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports; and

(b) means for connecting the switch nodes together in a multistage interconnect network, the means for connecting comprising forward channel and back channel signal paths coupled to each of the input and output ports in the switch nodes, wherein the back channel signal paths have a narrower bandwidth relative to the forward channel signal paths to simplify packaging.

2. The system of claim 1, further comprising:

(c) multicast means, operative within the network, for transmitting forward channel messages from a source to one or more destinations; and

(d) back channel merge means, within each switch node, for combining back channel replies received from the destinations into a single result, wherein the result is transmitted on the back channel to the source.

3. The system of claim 1, wherein the means for connecting comprises cabling means for wiring between switch nodes in different stages so that the forward channels and back channel signal paths are laid out in one or more copies of a universal wiring pattern.

4. The system of claim 3, wherein the universal wiring pattern comprises a permutation of connections between switch node input and output ports that swaps the least significant two base b digits of a port's level number representation, where b is the switch node size.

5. The system of claim 3, wherein the means for wiring comprises means for wiring every stage in a network of size $N=x^n$, $n>1$, with x^{n-2} copies of the universal wiring pattern, wherein x is a total number of output ports on a switch node.

6. The system of claim 1, wherein the means for connecting comprises means for interconnecting the switch nodes together to simplify the cabling, wherein a pattern of interconnections between switch nodes in each stage of the network is completely specified by permuting the digits of a level number representing a port of a switch node.

7. The system of claim 6, wherein the means for interconnecting comprises means for connecting switch node output ports identified by $(S : x_{n-1} \dots x_1 x_0)$ to switch node input ports identified by $(S+1 : \text{PERMUTE}_S^n \{x_{n-1} \dots x_1 x_0\})$, wherein S indicates a switch node in a specific stage, n refers to a total number of stages in the network, and x is a level number of a switch node port represented as $(x_{n-1} \dots x_1 x_0)_b$ in a base b corresponding to a size of the switch nodes, wherein $0 \leq x_i < b$, $0 \leq i < n$.

8. The network of claim 7, wherein $\text{PERMUTE}_S^n \{x_{n-1} \dots x_1 x_0\}$ is equivalent to $\text{PERMUTE}_0^2 \{x_1 x_0\} = x_0 x_1$ for a network wherein $n=2$.

9. The network of claim 7, wherein $\text{PERMUTE}_S^n \{x_{n-1} \dots x_1 x_0\}$ is equivalent to $\text{PERMUTE}_0^3 \{x_2 x_1 x_0\} = x_2 x_0 x_1$ and $\text{PERMUTE}_1^3 \{x_2 x_1 x_0\} = x_1 x_0 x_2$ for a network wherein $n=3$.

10. The network of claim 7, wherein $\text{PERMUTE}_S^n \{x_{n-1} \dots x_1 x_0\}$ is equivalent to $\text{PERMUTE}_0^3 \{x_2 x_1 x_0\} = x_2 x_0 x_1$ and $\text{PERMUTE}_1^3 \{x_2 x_1 x_0\} = x_0 x_1 x_2$ for a network wherein $n=3$.

11. The network of claim 7, wherein $\text{PERMUTE}_s^n \{x_{n-1} \dots x_1 x_0\}$ is equivalent to $\text{PERMUTE}_0^4 \{x_3 x_2 x_1 x_0\} = x_3 x_2 x_0 x_1$, $\text{PERMUTE}_1^4 \{x_3 x_2 x_1 x_0\} = x_1 x_0 x_3 x_2$, and $\text{PERMUTE}_2^4 \{x_3 x_2 x_1 x_0\} = x_3 x_2 x_0 x_1$ for a network wherein $n=4$.

12. The network of claim 11, wherein PERMUTE_0^4 and PERMUTE_2^4 leave the most significant two digits undisturbed, so that an interconnection length in stages 0 and 2 are minimized.

13. The system of claim 6, wherein the means for interconnecting comprises means for connecting switch node output ports identified by $(S+1 : \text{PERMUTE}_s^n \{x_{n-1} \dots x_1 x_0\})$ to switch node input ports identified by $(S : x_{n-1} \dots x_1 x_0)$, wherein S indicates a switch node in a specific stage, n refers to a total number of stages in the network, and x is a level number of a switch node port represented as $(x_{n-1} \dots x_1 x_0)_b$ in a base b corresponding to a size of the switch nodes, wherein $0 \leq x_i < b$ and $0 \leq i < n$.

14. The network of claim 13, wherein $\text{PERMUTE}_s^n \{x_{n-1} \dots x_1 x_0\}$ is equivalent to $\text{PERMUTE}_0^2 \{x_1 x_0\} = x_0 x_1$ for a network wherein $n=2$.

15. The network of claim 13, wherein
 $\text{PERMUTE}_s^n \{x_{n-1} \dots x_1 x_0\}$ is equivalent to $\text{PERMUTE}_0^3 \{x_2 x_1 x_0\} = x_2 x_0 x_1$
 and $\text{PERMUTE}_1^3 \{x_2 x_1 x_0\} = x_1 x_0 x_2$ for a network wherein $n=3$.

16. The network of claim 13, wherein
 $\text{PERMUTE}_s^n \{x_{n-1} \dots x_1 x_0\}$ is equivalent to $\text{PERMUTE}_0^3 \{x_2 x_1 x_0\} = x_2 x_0 x_1$
 and $\text{PERMUTE}_1^3 \{x_2 x_1 x_0\} = x_0 x_1 x_2$ for a network wherein $n=3$.

17. The network of claim 13, wherein
 $\text{PERMUTE}_s^n \{x_{n-1} \dots x_1 x_0\}$ is equivalent to $\text{PERMUTE}_0^4 \{x_3 x_2 x_1 x_0\} =$
 $x_3 x_2 x_0 x_1$, $\text{PERMUTE}_1^4 \{x_3 x_2 x_1 x_0\} = x_1 x_0 x_3 x_2$, and $\text{PERMUTE}_2^4 \{x_3 x_2 x_1 x_0\} =$
 $x_3 x_2 x_0 x_1$ for a network wherein $n=4$.

18. The network of claim 17, wherein PERMUTE_0^4 and
 PERMUTE_2^4 leave the most significant two digits undisturbed,
 so that an interconnection length in stages 0 and 2 are
 minimized.

19. The system of claim 1, wherein the means for connecting comprises means for arranging the switch nodes in more than $\lceil \log_b N \rceil$ stages, wherein b is a total number of switch node input/output ports, N is a total number of network I/O ports, and $\lceil \log_b N \rceil$ indicates a ceiling function providing the smallest integer not less than $\log_b N$, thereby providing additional communication paths between any network input port and network output port to enhance fault tolerance and lessen contention.

20. The system of claim 19, wherein the switch nodes are arranged in $2 \lceil \log_b N \rceil$ stages.

21. A method of communicating among a plurality of switch nodes connected together in a multistage interconnect network, each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports, the method comprising:

(a) transmitting messages between switch nodes using a forward channel coupled to each of the input and output ports in the switch nodes; and

(b) transmitting between switch nodes using a back channel coupled to each of the input and output ports in the switch nodes, wherein the back channels have a narrower bandwidth relative to the forward channels to simplify packaging.

22. The method of claim 21, further comprising:

(c) transmitting forward channel messages from a source to one or more destinations; and

(d) combining back channel replies received from the destinations into a single result, wherein the result is transmitted on the back channel to the source.

23. A communication system having a plurality of ports for concurrently transferring messages between different ones of the ports, comprising:

(a) a plurality of agents connected to the system, each including means providing interconnect request messages containing routing tags for communicating with other agents, means for transmitting variable length data messages, and means for transmitting response messages of different types having values in accordance with a predetermined protocol; and

(b) network means for intercoupling the agents, the network means comprising a plurality of switch nodes arranged in a multistage interconnect network and each having a plurality of forward channel and back channel signal paths with selectable interconnections between different terminals of the switch nodes, the back channel signal paths having a narrower bandwidth relative to the forward channel signal paths, the switch nodes also including means responsive to the routing tags in the messages for selecting signal paths between switch nodes in order to establish a complete path through the network, and the switch nodes further comprising means for reserving the complete path during the transmission of variable length data messages, once request messages and response messages thereto have been successfully interchanged.

24. A communications system, comprising:

(a) a plurality of switch nodes, each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports; and

(b) means for interconnecting the switch nodes together in a multistage interconnect network to simplify the cabling therebetween, wherein a pattern of interconnections between different stages of switch nodes is specified by permuting the digits of a level number representing a port of a switch node.

25. The system of claim 24, wherein the means for interconnecting comprises means for connecting switch node ports identified by $(S : x_{n-1} \dots x_1 x_0)$ to switch node ports identified by $(S+1 : \text{PERMUTE}_S^n \{x_{n-1} \dots x_1 x_0\})$, wherein S indicates a switch node in a specific stage, n refers to a total number of stages in the network, and x is a level number of a switch node port represented as $(x_{n-1} \dots x_1 x_0)_b$ in a base b corresponding to a size of the switch nodes, wherein $0 \leq x_i < b$ and $0 \leq i < n$.

26. The network of claim 25, wherein $\text{PERMUTE}_S^n \{x_{n-1} \dots x_1 x_0\}$ is equivalent to $\text{PERMUTE}_{x_0}^2 \{x_1 x_0\} = x_0 x_1$ for a network wherein $n=2$.

27. The network of claim 25, wherein $\text{PERMUTE}_s^a \{x_{n-1} \dots x_1 x_0\}$ is equivalent to $\text{PERMUTE}_0^3 \{x_2 x_1 x_0\} = x_2 x_0 x_1$ and $\text{PERMUTE}_1^3 \{x_2 x_1 x_0\} = x_1 x_0 x_2$ for a network wherein $n=3$.

28. The network of claim 25, wherein $\text{PERMUTE}_s^a \{x_{n-1} \dots x_1 x_0\}$ is equivalent to $\text{PERMUTE}_0^3 \{x_2 x_1 x_0\} = x_2 x_0 x_1$ and $\text{PERMUTE}_1^3 \{x_2 x_1 x_0\} = x_0 x_1 x_2$ for a network wherein $n=3$.

29. The network of claim 25, wherein $\text{PERMUTE}_s^a \{x_{n-1} \dots x_1 x_0\}$ is equivalent to $\text{PERMUTE}_0^4 \{x_3 x_2 x_1 x_0\} = x_3 x_2 x_0 x_1$, $\text{PERMUTE}_1^4 \{x_3 x_2 x_1 x_0\} = x_1 x_0 x_3 x_2$, and $\text{PERMUTE}_2^4 \{x_3 x_2 x_1 x_0\} = x_3 x_2 x_0 x_1$ for a network wherein $n=4$.

30. The network of claim 29, wherein PERMUTE_0^4 and PERMUTE_2^4 leave the most significant two digits undisturbed, so that an interconnection length in stages 0 and 2 are minimized.

31. A communications apparatus, comprising:

(a) a plurality of switch nodes, each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports; and

(b) means for connecting the switch nodes together in a multistage interconnect network by connecting switch node ports identified by $(S : x_{n-1} \dots x_1 x_0)$ to switch node ports identified by $(S+1 : \text{PERMUTE}_S^n \{x_{n-1} \dots x_1 x_0\})$, wherein S indicates a switch node in a specific stage, n refers to a total number of stages in the network and hence its size, and x is a level number of a switch node port represented as $(x_{n-1} \dots x_1 x_0)_b$ in a base b corresponding to the size of the switch nodes, wherein $0 \leq x_i < b$, $0 \leq i < n$ and $\text{PERMUTE}_S^n \{x_{n-1} \dots x_1 x_0\}$ is equivalent to $\text{PERMUTE}_0^2 \{x_1 x_0\} = x_0 x_1$ for a network wherein $n=2$.

32. A communications apparatus, comprising:

(a) a plurality of switch nodes, each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports; and

(b) means for connecting the switch nodes together in a multistage interconnect network by connecting switch node ports identified by $(S : x_{n-1} \dots x_1 x_0)$ to switch node ports identified by $(S+1 : \text{PERMUTE}_s^n \{x_{n-1} \dots x_1 x_0\})$, wherein S indicates a switch node in a specific stage, n refers to a total number of stages in the network and hence its size, and x is a level number of a switch node port represented as $(x_{n-1} \dots x_1 x_0)_b$ in a base b corresponding to the size of the switch nodes, wherein $0 \leq x_i < b$, $0 \leq i < n$ and $\text{PERMUTE}_s^n \{x_{n-1} \dots x_1 x_0\}$ is equivalent to $\text{PERMUTE}_0^3 \{x_2 x_1 x_0\} = x_2 x_0 x_1$ and $\text{PERMUTE}_1^3 \{x_2 x_1 x_0\} = x_1 x_0 x_2$ for a network wherein $n=3$.

33. A communications apparatus, comprising:

(a) a plurality of switch nodes, each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports; and

(b) means for connecting the switch nodes together in a multistage interconnect network by connecting switch node ports identified by $(S : x_{n-1} \dots x_1 x_0)$ to switch node ports identified by $(S+1 : \text{PERMUTE}_S^n \{x_{n-1} \dots x_1 x_0\})$, wherein S indicates a switch node in a specific stage, n refers to a total number of stages in the network and hence its size, and x is a level number of a switch node port represented as $(x_{n-1} \dots x_1 x_0)_b$ in a base b corresponding to the size of the switch nodes, wherein $0 \leq x_i < b$, $0 \leq i < n$ and $\text{PERMUTE}_S^n \{x_{n-1} \dots x_1 x_0\}$ is equivalent to $\text{PERMUTE}_0^3 \{x_2 x_1 x_0\} = x_2 x_0 x_1$ and $\text{PERMUTE}_1^3 \{x_2 x_1 x_0\} = x_0 x_1 x_2$ for a network wherein $n=3$.

34. A communications apparatus, comprising:

(a) a plurality of switch nodes, each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports; and

(b) means for connecting the switch nodes together in a multistage interconnect network by connecting switch node ports identified by $(S : x_{n-1} \dots x_1 x_0)$ to switch node ports identified by $(S+1 : \text{PERMUTE}_S^n \{x_{n-1} \dots x_1 x_0\})$, wherein S indicates a switch node in a specific stage, n refers to a total number of stages in the network and hence its size, and x is a level number of a switch node port represented as $(x_{n-1} \dots x_1 x_0)_b$ in a base b corresponding to the size of the switch nodes, wherein $0 \leq x_i < b$, $0 \leq i < n$ and $\text{PERMUTE}_S^n \{x_{n-1} \dots x_1 x_0\}$ is equivalent to $\text{PERMUTE}_0^4 \{x_3 x_2 x_1 x_0\} = x_3 x_2 x_0 x_1$, $\text{PERMUTE}_1^4 \{x_3 x_2 x_1 x_0\} = x_1 x_0 x_3 x_2$, and $\text{PERMUTE}_2^4 \{x_3 x_2 x_1 x_0\} = x_3 x_2 x_0 x_1$ for a network wherein $n=4$.

35. A communications network, comprising:

(a) a plurality of switch nodes arranged into a multistage interconnect network having a plurality of input and output ports, each port being coupled to an agent to effect communication between agents through the network;

(b) the network having a multiple of $\lceil \log_b N \rceil$ stages of interconnected switch nodes, wherein b is a total number of switch node input/output ports, N is the number of network I/O ports, and $\lceil \log_b N \rceil$ indicates a ceiling function providing the smallest integer not less than $\log_b N$, thereby providing additional paths between any network input port and network output port to enhance fault tolerance and lessen contention; and

(c) the network having a loop-back point indicating where the stages of the network are physically folded together so that corresponding switch nodes in similarly numbered stages on either side of the loop-back point are located adjacent to each other, thereby simplifying packaging and minimizing signal path lengths.

36. An apparatus for concurrently transferring messages between different ports, comprising:

(a) multistage interconnect network means for interconnecting a plurality of switch nodes for communication therebetween, each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports; and

(b) dynamic configuration means for determining how the switch nodes are interconnected and hence a topology of the multistage interconnect network means, so that messages can be routed correctly between the switch nodes.

37. The apparatus of claim 36, wherein the dynamic configuration means further comprises:

(1) means for communicating addresses between switch nodes, wherein one switch node communicates its location to another switch node connected thereto; and

(2) tag mapping table means, in each switch node, for storing routing information derived from the addresses of the switch nodes, so that messages can be routed correctly between the switch nodes.

38. The apparatus of claim 37, further comprising initializing means for querying the switch nodes at startup to determine how they are interconnected and for generating the tag mapping table means in response thereto.

39. A network, comprising:

(a) a plurality of switch nodes, each having a first plurality of input ports, a second plurality of output ports, and path selector means for selectively connecting the switch node input ports to the switch node output ports;

(b) means for interconnecting the switch nodes in a relatively arbitrary manner to effect a multistage interconnect network, the network having a first plurality of input ports, a second plurality of output ports, and means for routing messages through the network by transmitting a message from one switch node to another switch node; and

(c) means for determining how the switch nodes have been interconnected and for constructing routing tables for each switch node in response thereto, so that the path selector means can connect an input port receiving a message to an output port for transferring the message in order to correctly transmit the message.

40. The network of claim 39, wherein the means for determining comprises:

(1) topology determination means for communicating addresses between the switch nodes, so that a topology for the network can be determined; and

(2) message routing means for storing routing information derived from the addresses of the switch nodes, so that messages can be routed correctly through the network.

41. A communications apparatus, comprising:

(a) network means for providing bidirectional data transmission between agents, the network means comprising switch nodes connected together in a multistage interconnect network, each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports; and

(b) tag mapping table means, in each switch node, for storing routing information derived from the addresses of the switch nodes, so that the data transmission can be routed correctly through the network means.

42. A method for communicating in a network comprising a plurality of switch nodes, each switch node having a first plurality of input ports, a second plurality of output ports, and path selector means for selectively connecting the switch node input ports to the switch node output ports, the method comprising:

(a) interconnecting the switch nodes in a relatively arbitrary manner to effect a multistage interconnect network;

(b) determining how the switch nodes have been interconnected;

(c) constructing routing tables for each switch node according to how the switch nodes have been interconnected; and

(d) transferring messages through the network according to the routing tables, wherein each switch node that receives the messages uses a routing table to determine which output port should receive the message.

43. The method of claim 42, wherein the determining step comprises:

(1) communicating addresses between the switch nodes, so that a topology for the network can be determined; and

(2) storing routing information derived from the addresses of the switch nodes, so that messages can be routed correctly through the network.

44. A communications system, comprising:

(a) network means for providing bidirectional data transmission between agents connected thereto, the network means comprising switch nodes connected together in a multistage interconnect network, each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports; and

(b) self-diagnosing means, integrated with the network, for detecting and reporting any errors that occur within the network.

45. The system of claim 44, wherein the self-diagnosing means comprises diagnostic processor means for monitoring the state of the network means, for performing self-tests on the components of the network means, and for initializing the network means.

46. The system of claim 45, wherein the diagnostic processor means comprises means for configuring routing tables and input and output enable vectors for the network means so that fault conditions in the network means can be bypassed.

47. A communications system, comprising:

(a) a multistage interconnect network comprising a plurality of interconnected active logic switch nodes;

(b) diagnostic means for detecting and reporting any errors that occur within the network, and for isolating the errors without propagating them, thereby improving diagnosability and serviceability;

(c) reconfiguration means for reconfiguring the network when an error is detected, without interrupting communications in the system, so that any degradation in performance is minimized.

48. The system of claim 47, wherein the reconfiguration means further comprises means for reconfiguring routing tables and input and output enable vectors in the network, so that all communications are routed around a faulty section of the network.

49. The system of claim 47, wherein the network further comprises redundant networks for providing added bandwidth and fault tolerance to the system.

50. The system of claim 49, further comprising control means for load leveling message transmission through the redundant networks.

51. The system of claim 49, further comprising control means for switching between the redundant networks when a failure occurs, so that if one or more of the redundant networks is not available, the others can take over, thereby providing for graceful degradation in the presence of malfunctions.

52. A method for communicating in a multistage interconnect network comprising a plurality of interconnected active logic switch nodes, the method comprising the steps of:

(a) detecting and reporting any errors that occur within the network, and isolating the errors without propagating them, thereby improving serviceability; and

(b) reconfiguring the network when an error is detected, without interrupting communications in the system, so that any degradation in performance in the reconfigured network is minimized.

53. The method of claim 52, wherein the reconfiguring step further comprises reconfiguring routing tables and input and output enable vectors in the network, so that all communications are routed around a faulty section of the network.

54. The method of claim 52, further comprising load leveling message transmission through redundant networks.

55. The method of claim 52, further comprising switching between redundant networks when a failure occurs, so that if one or more of the redundant networks is not available, the others can take over, thereby providing for graceful degradation in the presence of malfunctions.

56. A communications system, comprising:

(a) a multistage interconnect network comprising a plurality of interconnected active logic switch nodes;

(b) each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports;

(c) the multistage interconnect network comprising more than $\lceil \log_b N \rceil$ stages of switch nodes, wherein b is a total number of switch node input/output ports, N is a total number of network input/output ports, and $\lceil \log_b N \rceil$ indicates a ceiling function providing the smallest integer not less than $\log_b N$, the stages thereby providing a plurality of paths between any network input port and network output port to enhance fault tolerance and lessen contention;

(d) diagnostic means for detecting and reporting any errors that occur within the network, thereby improving serviceability; and

(e) reconfiguration means for reconfiguring the network, without interrupting the communications in the system, when an error is detected, and for isolating the error without propagating it, so that any degradation in performance in the reconfigured network is minimized.

57. The system of claim 56, wherein the switch nodes further comprise tag mapping table means, associated with each input port of a switch node, for interpreting a routing tag to determine which output port of the switch node to select in order to route correctly a connect request through the network.

58. The system of claim 57, wherein the tag mapping table means comprises a memory array with a plurality of entries for translating the routing tag to an output port selection, wherein the array provides a one-to-one mapping between a logical port selection provided by the routing tag and a physical port selection.

59. The system of claim 58, wherein each entry in the array is derived from an appropriate field of the routing tag according to a stage of the switch node within the network.

60. The system of claim 56, wherein the switch nodes further comprise input enable vectors for indicating which input ports of the switch node are operational.

61. The system of claim 56, wherein the switch nodes further comprise output enable vectors for indicating which output ports of the switch node are operational.

62. The system of claim 59, wherein the tag mapping table means comprises a lookup table for mapping the routing tag to a physical output port based on the way in which the network is cabled.

63. The system of claim 59, wherein the tag mapping table means further comprises means for initializing each entry in the tag mapping table means to be numerically equivalent to its address offset, such that a physical port selection is equal to a logical port selection, thereby providing default values that are appropriate for fully configured networks.

64. The system of claim 63, wherein the tag mapping table means further comprises means for overlaying the default values when a configuring agent has determined an actual topology for the network.

65. A communications system, comprising:

(a) a network comprising at least two input ports and two output ports, and capable of simultaneous communications between a different pair of input and output ports;

(b) message routing means, within the network, for accepting a connect request containing a routing tag identifying an output port destination for the message, and for steering the connect request to the output port destination;

(c) back-off means, within the network, for cancelling the connect request when the output port destination is unavailable; and

(d) retry means, triggered by the back-off means, for delaying a retry of the connect request, and for trying a different connect request between the input port and a different output port destination, thereby reducing contention in the network.

66. The system of claim 65, wherein the network further comprises:

(1) a plurality of switch nodes, each comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports;

(2) tag mapping table means, associated with each input port of a switch node, for interpreting the routing tag to determine which output port of the switch node to select in order to route correctly the connect request through the network.

67. The system of claim 65, wherein the message routing means further comprises monocast blocking means for disabling the back-off means and for waiting until a desired path through the network becomes available.

68. The system of claim 65, wherein the message routing means further comprises monocast pipelining means for sending messages from sending agents to receiving agents without waiting for acknowledgments to connect requests from the receiving agents.

69. A communications method for a network comprising at least two input ports and two output ports, and capable of simultaneous communications between a different pair of input and output ports, the method comprising the steps of:

(a) accepting a connect request containing a routing tag identifying an output port destination for the message, and steering the connect request to the output port destination;

(b) cancelling the connect request when the output port destination is unavailable; and

(c) delaying a retry of the connect request, and for trying a different connect request from the input port, thereby reducing contention in the network.

70. The method of claim 69, further comprising interpreting the routing tag to determine which output port of the switch node to select in order to route correctly the connect request through the network.

71. The method of claim 69, wherein the accepting step further comprises preventing cancellation of connect requests and waiting until a desired path through the network becomes available.

72. The method of claim 69, wherein the accepting step further comprises transmitting messages from sending agents to receiving agents without waiting for acknowledgments to connect requests from the receiving agents.

RECEIVED
FEB 10 1964
U.S. DEPT. OF COMMERCE
BUREAU OF STANDARDS

73. A communications apparatus, comprising:

(a) a multistage interconnect network comprising a plurality of interconnected active logic switch nodes;

(b) each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports;

(c) the multistage interconnect network comprising more than $\lceil \log_b N \rceil$ stages of switch nodes, wherein b is a total number of switch node input/output ports, N is a total number of network input/output ports, and $\lceil \log_b N \rceil$ indicates a ceiling function providing the smallest integer not less than $\log_b N$, thereby providing a plurality of paths between any network input port and network output port to enhance fault tolerance and lessen contention; and

(d) load balancing means, in each switch node, for distributing messages among the plurality of output ports so that messages are evenly distributed throughout the network.

74. The apparatus of claim 73, wherein the load balancing means comprises:

(1) means, within each switch node, for choosing a switch node output port similarly numbered as a requesting switch node input port when the similarly numbered switch node output port is available; and

(2) means, within each switch node, for choosing a next available switch node output port when the similarly numbered switch node output port is not available.

75. The apparatus of claim 74, wherein the load balancing means further comprises means for selecting a second output port of a switch node having a next higher port number than the input port of the switch node that received the connect request when the first output port is unavailable.

76. A system for concurrently transferring messages between different ports, comprising:

(a) a plurality of switch nodes, each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting the input ports to the output ports; and

(b) means for interconnecting the switch nodes in a multistage interconnect network having a first plurality of network input ports and a second plurality of network output ports;

(c) partitioning means for grouping the network ports into logically independent subsets, wherein each subset is a supercluster; and

(d) multicast means, operative within the network, for transmitting a message from a network input port to one or more network output ports grouped in a supercluster, wherein messages transmitted within any one supercluster are prevented from interfering with messages transmitted within any other supercluster.

77. The system of claim 76, wherein the partitioning means further comprises:

(1) means for grouping network ports into a supercluster of size 2^{m-p} , wherein binary addresses for the grouped network ports are identical in p highest order bits thereof, thereby providing 2^p superclusters of size 2^{m-p} in the system, wherein $m = \lceil \log_2 N \rceil$, N is a total number of network input/output ports, and $\lceil \log_2 N \rceil$ indicates a ceiling function providing the smallest integer not less than $\log_2 N$; and

(2) means for addressing each supercluster as $\{Y_{m-1}, Y_{m-2}, \dots, Y_{m-p}\}$, wherein $y_i \in \{0, 1\}$ and $m-1 \leq i \leq m-p$.

78. The system of claim 77, further comprising means for recursively subdividing superclusters into smaller superclusters, so that the system can contain a plurality of superclusters of different sizes, wherein each size is a power of two.

79. The system of claim 77, further comprising means for allocating network ports to address blocks that are a power of two when each supercluster size is not a power of two and/or N is not a power of two.

80. The system of claim 76, wherein the partitioning means further comprises:

(1) means for generating a list of supercluster sizes;
and

(2) means for calculating a smallest power of two not less than each supercluster size;

(3) means for adding up the calculated powers of two to determine a size of a network needed;

(4) means for calculating the smallest power of two not less than the size of the network;

(5) means for dividing the network in half recursively as needed until there is a section which is the size of each power of two that was calculated for each supercluster; and

(6) means for assigning the network ports in each supercluster to addresses in the corresponding range.

81. A method for concurrently transferring messages between different ports of a multistage interconnect network having a plurality of interconnected switch nodes, the method comprising the steps of:

(a) grouping the network ports into logically independent subsets, wherein each subset is a supercluster; and

(b) transmitting a message from a network input port to one or more network output ports grouped in a supercluster, wherein messages transmitted within any one supercluster are prevented from interfering with messages transmitted within any other supercluster.

82. The method of claim 81, wherein the grouping step further comprises:

(1) grouping network ports into a supercluster of size 2^{m-p} , wherein binary addresses for the grouped network ports are identical in p highest order bits thereof, thereby providing 2^p superclusters of size 2^{m-p} in the system, wherein $m = \lceil \log_2 N \rceil$, N is a total number of network input/output ports, and $\lceil \log_b N \rceil$ indicates a ceiling function providing the smallest integer not less than $\log_b N$; and

(2) addressing each supercluster as $\{y_{m-1}, y_{m-2}, \dots, y_{m-p}\}$, wherein $y_i \in \{0, 1\}$ and $m-1 \leq i \leq m-p$.

83. The method of claim 82, further comprising recursively subdividing superclusters into smaller superclusters, so that the network can contain a plurality of superclusters of different sizes, wherein each size is a power of two.

84. The method of claim 82, further comprising allocating network ports to address blocks that are a power of two when N is not a power of two.

85. The method of claim 81, wherein the grouping step further comprises:

- (1) generating a list of supercluster sizes; and
- (2) calculating a smallest power of two not less than each supercluster size;
- (3) adding up the calculated powers of two to determine a size of a network needed;
- (4) calculating the smallest power of two not less than the size of the network;
- (5) dividing the network in half recursively as needed until there is a section which is the size of each power of two that was calculated for each supercluster; and
- (6) assigning the network ports in each supercluster to addresses in the corresponding range.

86. A system for concurrently transferring messages between different ports, comprising:

(a) a plurality of switch nodes, each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports; and

(b) means for connecting the switch nodes together in a multistage interconnect network, the means for connecting comprising forward channel and back channel signal paths; and

(c) multicast means, operative within the network, for transmitting forward channel messages from a source to one or more destinations; and

(d) back channel merge means, within each switch node, for combining back channel replies received from the destinations into a single result, wherein the result is transmitted on the back channel to the source.

87. The system of claim 86, wherein the multicast means comprises means for steering a multicast request for a supercluster to a bounce back point within the network means, wherein all multicast requests to the supercluster use the same bounce back point.

88. The system of claim 87, wherein the means for steering comprises means for steering a multicast request from one supercluster to a destination supercluster through a bounce back point for the destination supercluster.

89. The system of claim 86, wherein the multicast means comprises means for permitting only one multicast message at a time within a supercluster, thereby preventing deadlock between competing multicast requests.

90. A method for concurrently transferring messages between different ports of multistage interconnect network, the network comprising a plurality of switch nodes, each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports, the switch nodes connected together via forward channel and back channel signal paths connected to every input and output port, the method comprising the steps of:

(a) transmitting forward channel messages from a source to one or more destinations; and

(b) combining back channel replies received from the destinations into a single result, wherein the result is transmitted on the back channel to the source.

91. The method of claim 90, wherein the transmitting step comprises steering a multicast request for a supercluster to a bounce back point within the network means, wherein all multicast requests to the supercluster use the same bounce back point.

92. The method of claim 91, wherein the steering step comprises steering a multicast request from one supercluster to a destination supercluster through a bounce back point for the destination supercluster.

94. A system for concurrently transferring messages, comprising:

(a) a multistage interconnect network comprising a plurality of interconnected active logic switch nodes;

(b) each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports;

(c) the multistage interconnect network comprising more than $\lceil \log_b N \rceil$ stages of switch nodes, wherein b is a total number of switch node input/output ports, N is a total number of network input/output ports, and $\lceil \log_b N \rceil$ indicates a ceiling function providing the smallest integer not less than $\log_b N$, the multistage interconnect network providing a plurality of paths between any network input port and network output port to enhance fault tolerance and lessen contention; and

(d) multicast steering means, within each switch node, for routing multicast requests to a specific input port of a specific switch node within the network, so that only one multicast request can occur at a time, thereby preventing deadlock between competing multicast requests.

95. The system of claim 94, further comprising:

(1) means for storing a reply from each network output port in the back channel; and

(2) means for collecting replies from the network output ports and for applying the replies to merge means for synchronously combining all of the replies, wherein the replies are sorted as they propagate through the merge means, so that only the reply having the highest priority is transmitted through the system.

96. The system of claim 95, wherein the merge means comprises low sort means for outputting a reply with a lowest key followed by an accompanying data word.

97. The system of claim 95, wherein the merge means further comprises add means for summing data words of all replies in a bit serial manner so that the result has the same number of bits as the operands.

98. An apparatus for concurrently transferring messages between different ports, comprising:

(a) a multistage interconnect network comprising a plurality of interconnected active logic switch nodes;

(b) each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports,

(c) the multistage interconnect network comprising more than $\lceil \log_b N \rceil$ stages of switch nodes, wherein b is a total number of switch node input/output ports, N is a total number of network input ports and network output ports, and $\lceil \log_b N \rceil$ indicates a ceiling function providing the smallest integer not less than $\log_b N$, the multistage interconnect network providing a plurality of paths between any network input port and network output port to enhance fault tolerance and lessen contention; and

(d) deadlock avoidance means, within each switch node, for allowing only one routing multicast request at a time, thereby preventing deadlock between requests.

99. The system of claim 98, wherein the deadlock avoidance means further comprises means for allowing only one multicast request within a supercluster at a time and for allowing a plurality of multicast requests within different superclusters to occur at the same time.

100. A communications system, comprising:

(a) a plurality of switch nodes arranged into a multistage interconnect network having a plurality of input and output ports, each port being coupled to an agent to effect communication between agents through the network;

(b) the network having more than $\lceil \log_b N \rceil$ stages of interconnected switch nodes, wherein b is a total number of switch node input/output ports, N is the number of network input/output ports, and $\lceil \log_b N \rceil$ indicates a ceiling function providing the smallest integer not less than $\log_b N$; and

(c) the network having a plurality of turnaround points at the highest stage of switch nodes, the turnaround points logically differentiating between switch nodes that load balance messages through the network from switch nodes that direct messages to receiving agents;

(d) means for depopulating switch nodes from the highest stage to reduce the number of turnaround points in the network, as long as at least one path exists between every network input port and every network output port; and

(e) the input and output ports of the switch nodes in stages adjacent the highest stage sensing when the switch nodes in the highest stage are removed and disabling the input and output ports in response thereto, thereby lowering the bandwidth of the network and lowering the cost of the network without a loss of functionality.

101. A communications apparatus, comprising:

(a) multistage interconnect network means for interconnecting a plurality of switch nodes for communication therebetween;

(b) dynamic configuration means for determining how the switch nodes are interconnected by the multistage interconnect network means; and

(c) means for remapping connections between the switch nodes so that all connections in a backplane connecting the switch nodes are horizontal.

102. A network, comprising:

(a) a plurality of switch nodes, each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports;

(b)' backplane means for interconnecting the input and output ports of different switch nodes to effect a multistage interconnect network; and

(c) tag mapping table means for remapping the interconnections between the switch nodes so that all connections in the backplane means are horizontal only.

103. A communications systems, comprising:

(a) a network comprising a plurality of interconnected switch nodes; and

(b) identification means for communicating a unique identifier to each agent attached to the network, so that an agent can be connected to any available port of the network and determine its address therein, thereby simplifying installation.

104. The system of claim 103, wherein the identification means comprises means for determining the agent's address within the network means using a level number associated with the port to which it is connected.

105. The system of claim 104, wherein the means for determining comprises means for transmitting a command to the network means, wherein a switch node that receives the command replies thereto with a port address.

106. The system of claim 104, wherein the means for determining comprises response means, within a switch node connected to the input port, for appending a port number of the switch node to the level number of the switch node to create the agent address in the network means.

107. A communications system, comprising:

(a) network means for providing bidirectional data transmission between a plurality of agents connected thereto, the network means comprising switch nodes connected together in a multistage interconnect network, each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports; and

(b) means for arranging components within the network means and interconnections between switch nodes, so that the number of agents can be increased with only a substantially linear increase in the size of the network means.

108. The system of claim 107, wherein the network is folded and the switch nodes are paired so that corresponding switch nodes in a specific stage are physically adjacent to one another on a board.

109. The system of claim 107, further comprising means for maintaining physical interconnections between switch nodes at a minimal length.

110. The system of claim 107, wherein the means for arranging comprises:

(1) first means for constructing the network means using only a first board type when the network means contains 8 or fewer network I/O ports;

(2) second means for constructing the network means using only a first board type and a second board type when the network means contains between 9 and 64 network I/O ports;

(3) third means for constructing the network means using only a first board type, a second board type, and a third board type when the network means contains between 65 and 512 network I/O ports; and

(4) fourth means for constructing the network means using only a first board type, a second board type, and a fourth board type when the network means contains between 513 and 4096 network I/O ports.

111. A communications system, comprising:

(a) network means, coupled to the agents, for providing bidirectional data transmission between a plurality of agents connected thereto, the network means comprising switch nodes connected together in a multistage interconnect network, each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports; and

(b) cabling means for wiring between different stages in the network means with one or more copies of a universal wiring pattern.

112. The system of claim 111, wherein the universal wiring pattern comprises a permutation of switch node ports that swaps the least significant two base b digits of a level number representation, where b is a total number of switch node input or output ports.

113. The system of claim 111, wherein the means for wiring comprises means for wiring every stage in a network means of size $N=x^n$, $n>1$, with x^{n-2} copies of the universal wiring pattern, wherein x is a number of output ports on a switch node.

114. A communications system, comprising:

(a) network means for providing bidirectional data transmission between a plurality of agents connected thereto, the network means comprising switch nodes connected together in a multistage interconnect network, each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports;

(b) master clock means for supplying a master clock signal to the switch nodes throughout the network means; and

(c) controller means, coupled between the network means and the agents, for communicating messages between the agent and the network means, the controller means comprising clock extraction means for deriving the master clock signal transmitted throughout the network means, wherein the master clock signal increments a counter in the controller means to provide a synchronous date and time of day to each agent.

115. A system for communicating in a variety of modes between one or more transmitting agents and one or more receiving agents, comprising:

(a) a multistage interconnect network intercoupling all the agents with transmit and receive lines, the network comprising a plurality of switch nodes arranged in parallel groupings in a plurality of stages within the network;

(b) the transmitting agents including means for transmitting message routing packets containing destination data designating one or more receiving agents and means for transmitting variable length data messages; and

(c) wherein the switch nodes each comprise means responsive to the message routing packets for selecting node-to-node paths to one or more receiving agents, wherein the switch nodes also comprise means for acknowledging receipt of a message routing packet and means for preserving the circuit connection for subsequent transmission of a variable length message.

116. An active logic network, comprising:

(a) a plurality of switch nodes, each having a plurality of "left-hand" ports and a plurality of "right-hand" ports, wherein within the switch node, any left hand port can be connected to any right hand port, and any right hand port can be connected to any left hand port;

(b) means, within each switch node, for accepting a routing tag at an input port identifying a destination in the network; and

(c) routing determination means, within each switch node, for translating the routing tag into an output port selection of the switch node, thereby identifying a path to the destination.

117. The network of claim 116, wherein the switch node further comprises:

(1) means for accepting a connect request designating a desired destination to which a connection is requested;

(2) selector means, connected to the means for accepting, for selecting an output port of the switch node according to the desired destination; and

(3) arbiter means, connected to the port selector means and associated with each output port, for transmitting the connect request to the selected output port, the arbiter means choosing between contending connect requests when more than one port selector means is trying to access the selected output port at the same time, wherein the connect request is transmitted to the selected output port when it has been granted priority over the contending connect requests.

118. The network of claim 117, wherein the selector means further comprises means for selecting an output port based on a table lookup.

119. The network of claim 117, wherein the selector means further comprises means for evenly distributing connect requests between a plurality of arbiter means.

120. The network of claim 117, wherein the arbiter means further comprises:

(1) means for servicing contending point-to-point connect requests on a "round robin" basis; and

(2) means for providing priority to multicast connect requests over point-to-point connect requests.

121. The network of claim 116, wherein the switch nodes further comprises:

(1) means for transmitting messages through a forward channel; and

(2) means for receiving responses through a back channel, wherein the back channel has a narrower bandwidth than the forward channel.

122. The network of claim 121, wherein the switch nodes further comprise:

(3) means for collecting responses in the back channel;

(4) means for synchronously combining the collected responses so that the collected responses are sorted as they propagate through the switch nodes, wherein only a response having a highest priority is transmitted through the back channel.

123. A communications system, comprising:

(a) network means, coupled to the agents, for providing bidirectional data transmission between network ports, the network means comprising switch nodes connected together in a multistage interconnect network; and

(b) arbiter means, in each switch node, for choosing between contending connect requests accepted concurrently from a plurality of switch node ports, wherein the contending connect requests are all trying to access a selected output port at the same time, so that a connect request is transmitted to the selected output port when it has been granted priority over the contending connect requests.

124. The system of claim 123, wherein the arbiter means further comprises means for evenly distributing requests between a plurality of arbiter means.

125. The system of claim 124, wherein the arbiter means further comprises:

(1) means for providing a "round robin" priority to the contending point-to-point connect requests; and

(2) means for providing priority to multicast connect requests over point-to-point connect requests.

126. A system for transmitting messages between agents in (1) an arbitrary interconnection mode or (2) a multicast mode, comprising:

(a) a multistage interconnect network comprising a plurality of multiple terminal bidirectional switch nodes arrayed in a plurality of stages;

(b) the agents each including means for generating addressing messages containing destination data, including alternative descriptors for individual and multicast group designations, and each being coupled to the switch nodes; and

(c) wherein the switch nodes further include means responsive to the descriptors in the addressing messages for selecting routing paths through the network and means responsive to path selection for establishing path commitments linking agents for communication of variable length messages.

127. A network for communicating between agents connected thereto, comprising:

(a) packet switching means for establishing a communication path between sending and receiving agents in response to a connect request;

(b) circuit switching means for transferring messages of arbitrary length between the sending and receiving processors once the communication path has been established;

(c) pipelining means for transferring messages between sending agents and the receiving agents without waiting for the communication path to be established; and

(d) back-off means for cancelling the connect request when the communication path to the receiving agent is unavailable; and

(e) retry means, triggered by the back-off means, for delaying a retry of the connect request, and for trying a different connect request from the sending agent, thereby reducing contention in the network.

128. The network of claim 127, wherein the packet switching means further comprises means for transferring a connect command between switch nodes in the communication path.

129. The network of claim 128, wherein links between switch nodes comprising the communication path are reserved until released by the sending and receiving agents.

130. A communications apparatus, comprising:

(a) a plurality of switch nodes arranged and interconnected into a multistage interconnect network; and

(b) error testing means for generating test patterns during data transmission through the network, the error testing means comprising means, within each switch node, for inverting parity on the data transmission so that a receiving switch node reports an error.

131. The apparatus of claim 130, wherein the error testing means further comprises means for forwarding a test command to the receiving switch node and means for inverting parity on data which follows the test command.

132. The apparatus of claim 130, wherein the error testing means further comprises means for forwarding a test reply to the receiving switch node with its parity inverted.

133. The apparatus of claim 130, wherein the error testing means further comprises means for inverting parity continuously on the data transmission to the receiving switch node.

134. A switch node providing for the simultaneous interconnection of multiple messages, comprising:

(a) means defining a plurality of input and output ports;
(b) means for selectively connecting any input port to (1) any one output port, (2) a predetermined number of the output ports, or (3) all of the output ports; and

(c) error testing means for generating test patterns to test data transmission on the input and output ports, the error testing means comprising means for inverting parity on the data transmission so that a switch node receiving the data transmission reports an error.

135. The switch node of claim 134, wherein the error testing means further comprises means for forwarding a test command to the receiving switch node and means for inverting parity on data which follows the test command.

136. The switch node of claim 134, wherein the error testing means further comprises means for forwarding a test reply to the receiving switch node with its parity inverted.

137. The switch node of claim 134, wherein the error testing means further comprises means for inverting parity continuously on the data transmission to the receiving switch node.

138. A switch node for use in a network having a plurality of similar nodes interconnected by cabling, for simultaneous interconnection of multiple messages of different types, including message routing packets, comprising:

(a) means defining a plurality of input and output terminals;

(b) means for storing reconfigurable mapping data which identifies available interconnection paths using the input and output terminals; and

(c) means, coupled to the means for storing, for arbitrarily interconnecting any input with (1) any one output, (2) a predetermined number of the outputs or (3) all the outputs.

139. A switch node as set forth in claim 138 above, wherein means for storing reconfigurable mapping data further comprises means for storing input and output enable vectors indicating which terminals are operation.

140. A switch node as set forth in claim 138 above, wherein the switch nodes each include forward channel and back channel signal paths coupled to each of the input and output ports in each of the switch nodes, wherein the back channel signal paths have a narrower bandwidth relative to the forward channel signal paths to simplify packaging.

141. A switch node as set forth in claim 140 above, further comprising means for combining signals on the different back channel signal paths within the switch node.

142. A switch node as set forth in claim 140 above, further comprising means for synchronously combining replies from the back channels so that replies are sorted as they propagate through the switch node, wherein only the reply having a highest priority is transmitted to the next switch node on the back channel.

143. A switch node as set forth in claim 138 above, wherein the means interconnecting the inputs and outputs include means for arbitrating between competing message routing packets to determine priority in accordance with predetermined rules.

144. A switch node as set forth in claim 143 above, wherein the means for arbitrating affords priority to messages whose routing packets indicate more than one receiving destination.

145. A switch node as set forth in claim 138 above, wherein said switch node includes means for leveling the load of messages through the switch node.

146. A switch node as set forth in claim 138 above, wherein the switch node includes means for returning a signal back to an input as to the unavailability of an interconnection.

147. A switch node as set forth in claim 138 above, wherein the switch node includes means for rerouting a message between different input and output terminals when a given path is unavailable.

148. A computer system, comprising:

(a) a plurality of processors comprising (1) application processors (AP) for connecting external devices to the system, and (2) access module processors (AMP) for managing access to a relational database; and

(b) network means, coupled to the processors, for providing bidirectional data transmission between the processors, the network means comprising switch nodes connected together in a multistage interconnect network, each switch node comprising a first plurality of input ports, a second plurality of output ports, and means for selectively connecting said input ports to said output ports, the network having more than $\lceil \log_b N \rceil$ stages of interconnected switch nodes, wherein b is a total number of switch node input/output ports, N is the number of network input/output ports, and $\lceil \log_b N \rceil$ indicates a ceiling function providing the smallest integer not less than $\log_b N$.

149. The computer system of claim 148, wherein the AP further comprises one or more microprocessors, including memory and at least two connections to the network, for executing an independent copy of the operating system and maintaining an independent address space.

150. The computer system of claim 149, wherein the AP further comprises a plurality of tightly-coupled microprocessors sharing a single copy of the operating system and sharing a common address space.

151. The computer system of claim 148, where the AMP further comprises one or more microprocessors having disk I/O capabilities and at least two connections to the network, and executing an independent copy of the operating system that is specifically designed for executing database software.

152. The computer system of claim 151, wherein the AMP further comprises a plurality of tightly-coupled microprocessors sharing a single copy of the operating system and sharing a common address space.

153. The computer system of claim 148, where the system further comprises means for allocating a portion of the relational database to each AMP.